

```
In [1]: import pandas as pd
import numpy as np
water=pd.read_csv("water_potability.csv")
type(water)
```

Out[1]:pandas.core.frame.DataFrame

```
In [2]: water.head()
```

Out[2]:

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
0	NaN	204.890455	20791.318981	7.300212	368.516441	564.308654	10.379783	86.990970	2.963135	0
1	3.716080	129.422921	18630.057858	6.635246	NaN	592.885359	15.180013	56.329076	4.500656	0
2	8.099124	224.236259	19909.541732	9.275884	NaN	418.606213	16.868637	66.420093	3.055934	0
3	8.316766	214.373394	22018.417441	8.059332	356.886136	363.266516	18.436524	100.341674	4.628771	0
4	9.092223	181.101509	17978.986339	6.546600	310.135738	398.410813	11.558279	31.997993	4.075075	0

```
In [5]: water.isnull().sum()
```

Out[5]:

ph	491
Hardness	0
Solids	0
Chloramines	0
Sulfate	781
Conductivity	0
Organic_carbon	0
Trihalomethanes	162
Turbidity	0
Potability	0
dtype:	int64

```
In [6]: water.describe()
```

Out[6]:

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
count	2785.000000	3276.000000	3276.000000	3276.000000	2495.000000	3276.000000	3276.000000	3114.000000	3276.000000	3276.000000
mean	7.080795	196.369496	22014.092526	7.122277	333.775777	426.205111	14.284970	66.396293	3.966786	0.390110
std	1.594320	32.879761	8768.570828	1.583085	41.416840	80.824064	3.308162	16.175008	0.780382	0.487849
min	0.000000	47.432000	320.942611	0.352000	129.000000	181.483754	2.200000	0.738000	1.450000	0.000000
25%	6.093092	176.850538	15666.690297	6.127421	307.699498	365.734414	12.065801	55.844536	3.439711	0.000000
50%	7.036752	196.967627	20927.833607	7.130299	333.073546	421.884968	14.218338	66.622485	3.955028	0.000000
75%	8.062066	216.667456	27332.762127	8.114887	359.950170	481.792304	16.557652	77.337473	4.500320	1.000000
max	14.000000	323.124000	61227.196008	13.127000	481.030642	753.342620	28.300000	124.000000	6.739000	1.000000

```
In [7]: water.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3276 entries, 0 to 3275
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0  ph                    2785 non-null  float64
1  Hardness              3276 non-null  float64
2  Solids                3276 non-null  float64
3  Chloramines          3276 non-null  float64
4  Sulfate               2495 non-null  float64
5  Conductivity          3276 non-null  float64
6  Organic_carbon        3276 non-null  float64
7  Trihalomethanes       3114 non-null  float64
8  Turbidity             3276 non-null  float64
9  Potability            3276 non-null  int64
dtypes: float64(9), int64(1)
memory usage: 256.1 KB
```

```
In [8]: water.dropna()
```

Out[8]:

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
3	8.316766	214.373394	22018.417441	8.059332	356.886136	363.266516	18.436524	100.341674	4.628771	0
4	9.092223	181.101509	17978.986339	6.546600	310.135738	398.410813	11.558279	31.997993	4.075075	0
5	5.584087	188.313324	28748.687739	7.544869	326.678363	280.467916	8.399735	54.917862	2.559708	0
6	10.223862	248.071735	28749.716544	7.513408	393.663396	283.651634	13.789695	84.603556	2.672989	0
7	8.635849	203.361523	13672.091764	4.563009	303.309771	474.607645	12.363817	62.798309	4.401425	0
...
3267	8.989900	215.047358	15921.412018	6.297312	312.931022	390.410231	9.899115	55.069304	4.613843	1
3268	6.702547	207.321086	17246.920347	7.708117	304.510230	329.266002	16.217303	28.878601	3.442983	1
3269	11.491011	94.812545	37188.826022	9.263166	258.930600	439.893618	16.172755	41.558501	4.369264	1
3270	6.069616	186.659040	26138.780191	7.747547	345.700257	415.886955	12.067620	60.419921	3.669712	1
3271	4.668102	193.681735	47580.991603	7.166639	359.948574	526.424171	13.894419	66.687695	4.435821	1

2011 rows × 10 columns

QUESTION 2 An outlier is a data point that significantly deviates from the majority of the data in a dataset. Here is an example: As a teacher, you recorded the exam scores (out of 100) of 5 students and the followings are: 85, 92, 77, 75, 13

In this scenario, "13" is the outlier because it is significantly lower than the rest of the scores.

Here are the importance of an outlier in data analysis: While using data visualization methods, outliers are easily spotted and can affect how the model would look like, sometimes it would also affect how trends are represented on a line chart.

Many statistical methods, such as linear regression, assume that the data follows a certain distribution. Outliers can misguide these assumptions leading to innacurate or biased results.

Identifying them can help signal errors in data collection, entry or processing. They can also reveal rare events, anomalies that are valuanle in some analysis.

QUESTION 3 The first common method in identifying an outlier in a dataset is to use visual representations. Scatter plots, Box plots and histograms are the ones that will make the outliers visible.

The next common one would be using the Z-score method. We can establish the lower and upper bounds of the dataset to discover outliers using the z score.

And the 3rd common one is called the Interquartile range,It is calculated by taking the difference between theupper and lower quartiles of the dataset. We can define the upper and lowerboundaries of a dataset using IQR. $Q3 + 1.5IQR$ is defined as the upper bound, and $Q1 - 1.5IQR$ is defined as the lower bound. Outliers are any observations or data points that fall outside of these parameters.

In []:

Loading [MathJax]/extensions/Safe.js